

# Exploring the Potential of an Extensible Domain-Specific Web Corpus for “Layfication”: The Case of Cross-Lingual Classification

Marina Santini, RISE Research Institutes of Sweden, Sweden

Min-Chun Shih, Linköping University, Sweden

## ABSTRACT

This article presents experiments based on the extensible domain-specific web corpus for “layfication”. For these experiments, both the existing layfication corpus (in Swedish and in English) and a new addition in English (the NHS-PubMed subcorpus) are used. With this extended corpus, methods to classify lay-specialized medical sublanguages cross-linguistically using small data and noisy web documents are investigated. Sublanguage is a language variety used in specific domains. Here, the authors focus on two medical sublanguages, namely the “patientspeak” (lay) and the medical jargon (specialized). Cross-lingual sublanguage classification is still largely underexplored although it can be crucial in downstream applications for digital health and cyber-physical systems. Classification models are built using small and noisy training sets in Swedish and evaluated on English test sets. The performance of Naive Bayes classifiers—built with stopwords and with Bag-of-Words—is compared with convolutional neural network classifiers leveraging on MUSE multi-lingual word embeddings. Results are promising and nuanced. These results are proposed as a first baseline for cross-lingual sublanguage classification.

## KEYWORDS

CNN, Convolutional Neural Network, Medical Domain, Naïve Bayes, Noise, Noisy Data, Small Data, Specialized Corpus, Sublanguage, Web Corpus

DOI: 10.4018/IJCPS.2020010102

Copyright © 2020, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

## INTRODUCTION

Cyber-Physical Systems (CPSs) belong to an emergent paradigm that combines most advanced technological approaches and computational tools to solve complex tasks. CPSs are domain-independent and have penetrated diversified disciplines, such as healthcare and self-driving vehicles.

In the era of data-driven science, corpus-based Language Technology is an essential component of many CPSs, where linguistic knowledge is indispensable to prevent failures or fatal errors due to misunderstandings or poor understanding during human communication, especially when it comes to the interaction between doctors and patients. Arguably, Language Technology can help because it empowers patients and other non-professional actors to understand medical information. Since modern Language Technology and Artificial Intelligence are data-driven, the need of representative text corpora is paramount. Curated, big and “clean” corpora are expensive and beyond the budget of practical domain-specific medical applications. Rather, real-world data is often small and noisy but nonetheless potentially useful. To address this issue, a textual resource was created in the form of an extensible web corpus for “layfication” (Santini et al. 2019) whose main purpose is to explore and test methods based on small and noisy data. The aim is to use the extensible web corpus for layfication as a sandbox to explore the benefits and the limitations of computational methods that could enhance digital health and Cyber-Physical Systems.

In Santini et al. (2019), a corpus for layfication was been presented and made available to the community for re-use, enhancement and expansion. The corpus was conceived as a flexible and extensible textual resource, where additional documents and additional languages can be appended over time. The corpus was started with web pages collected from the internet and includes documents describing a number of diseases in lay and specialized sublanguages. Although progressively expanded, the corpus is still small, because information about some diseases is not extensive. This paucity of textual data in some areas is a real-world condition. The main purpose of the corpus is to be used for building and training language technology applications for the “layfication” of the specialized medical jargon, where “layfication” refers to the automatic identification of more intuitive linguistic expressions that can help laypeople (e.g., patients, family caregivers, and home care aides) understand medical terms, which often appear opaque. In the medical field, layfication is also denoted by the expression “patientspeak”, to indicate that the layfication is needed by many patients in their interaction with healthcare professionals.

Why is layfication needed? The medical domain hinges upon medical terminology elaborated and used by healthcare professionals. Despite it is widely acknowledged that understanding what the doctors say has an important influence on the success of treatments, in many cases medical terminology hinders the comprehension of various groups of people (such as non-native speakers, people with low education, etc.), and has negative effects on health consumers (e.g. patients and caregivers). As a matter of fact, in the medical domain, two broad user groups interact. The first group (the expert)

11 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: [www.igi-global.com/article/exploring-the-potential-of-an-extensible-domain-specific-web-corpus-for-layfication/272559](http://www.igi-global.com/article/exploring-the-potential-of-an-extensible-domain-specific-web-corpus-for-layfication/272559)

## Related Content

---

### Addressing Special Educational Needs in Classroom With Cyber Physical Systems

Aneta Petrova Atanasova and Aleksandra Ivaylova Yosifova (2019). *Cyber-Physical Systems for Social Applications* (pp. 22-43).

[www.irma-international.org/chapter/addressing-special-educational-needs-in-classroom-with-cyber-physical-systems/224414](http://www.irma-international.org/chapter/addressing-special-educational-needs-in-classroom-with-cyber-physical-systems/224414)

### Semiotic Brains and Artificial Minds: How Brains Make up Material Cognitive Systems

Lorenzo Magnani (2007). *Semiotics and Intelligent Systems Development* (pp. 1-41).

[www.irma-international.org/chapter/semiotic-brains-artificial-minds/28935](http://www.irma-international.org/chapter/semiotic-brains-artificial-minds/28935)

### From Intermediary to Mediator and Vice Versa: On Agency and Intentionality of a Mundane Sociotechnical System

Antonio Díaz Andrade (2010). *International Journal of Actor-Network Theory and Technological Innovation* (pp. 21-29).

[www.irma-international.org/article/intermediary-mediator-vice-versa/47531](http://www.irma-international.org/article/intermediary-mediator-vice-versa/47531)

### Where is the Missing Matter?: A Comment on "The Essence" of Additive Manufacturing

Tihomir Mitev (2015). *International Journal of Actor-Network Theory and Technological Innovation* (pp. 10-17).

[www.irma-international.org/article/where-is-the-missing-matter/126277](http://www.irma-international.org/article/where-is-the-missing-matter/126277)

### The Structure of Theory and the Structure of Scientific Revolutions: What Constitutes an Advance in Theory?

Steven E. Wallis (2010). *Cybernetics and Systems Theory in Management: Tools, Views, and Advancements* (pp. 151-175).

[www.irma-international.org/chapter/structure-theory-structure-scientific-revolutions/39327](http://www.irma-international.org/chapter/structure-theory-structure-scientific-revolutions/39327)